

# Clasificación automática de evoluciones médicas multiclasas en español

Marcela Riccillo<sup>1</sup>, David Perez<sup>1</sup>, Daniel Luna<sup>1</sup>, Fernando Campos<sup>1</sup>,  
Carlos Otero<sup>1</sup>, María Laura Gambarte<sup>1</sup>, Sonia Benítez<sup>1</sup>

<sup>1</sup>Departamento de Informática en Salud, Hospital Italiano de Buenos Aires HIBA  
mail: marcela.riccillo@hospitalitaliano.org.ar

**Abstract.** En este trabajo presentamos una comparación de metodologías de clasificación para texto libre de narrativas médicas, en este caso evoluciones médicas multiclasas. Comparamos el rendimiento de redes neuronales y máquinas de soporte vectorial con preprocesamientos para clasificar evoluciones de Diabetes, en Tipo 1, Tipo 2 y Otros (otro tipo de afección). Se compararon accuracy, sensitivity y specificity, mostrando beneficios en costos de entrenamiento y resultados de exactitud. Encontramos porcentajes mayores con redes neuronales sin preprocesamiento PCA y en el caso de SVM con dicho preprocesamiento (con menor costo de entrenamiento).

**Keywords:** Electronic Health Record, Support Vector Machine, Natural Language Processing, Neural Network, Diabetes Mellitus, Narrative Medicine, Principal Component Analyses

## 1 Introducción

En el Hospital Italiano de Buenos Aires [1], desde el año 1998 se ha implementado de manera gradual un Sistema de Información en Salud (SIS) a partir de un desarrollo “in house” que maneja la información médica y administrativa desde la captura hasta el análisis. Incluye una única Historia Clínica Electrónica (HCE) web, modular, orientada a problemas y centrada en el paciente. Conocida con el nombre de ITALICA, la HCE permite el registro de la atención en los ámbitos: ambulatorio, internación, emergencias y atención domiciliaria. ITALICA permite la solicitud de estudios complementarios, prescripción farmacológica y visualización de resultados que incluye un sistema de almacenamiento y transmisión de imágenes asociadas al paciente.

Dado que la HCE es orientada a problemas, cuenta con una lista estructurada de los problemas de los pacientes que posee un control terminológico [2] que normaliza el texto ingresado según ontologías de referencia. Luego, el médico evoluciona la situación de cada paciente asociando dicha evolución al problema correspondiente. La evolución es un campo que contiene textos libres de narrativa médica.

En la narrativa médica encontramos información valiosa sobre las características de cada paciente, evolución, medicamentos, antecedentes familiares, entre otros, que podría ser aprovechada para complementar y en algunos casos hasta para cotejar los

datos que aparecen en forma estructurada. Esto comprende por ejemplo [3] informes radiológicos, resúmenes de alta, reportes de patologías, historias de admisión, informes de exámenes físicos.

Los textos libres en general son complejos con datos de laboratorio (donde algunos son copiados de otros campos para enfatizar o ilustrar las indicaciones), comportamientos del paciente, sugerencias de tratamientos, interconsultas, entre otros. La complejidad lingüística dificulta la extracción automática de información, que luego podría utilizarse para reportes estadísticos o la correcta determinación de la patología presentada por el paciente.

El Procesamiento del Lenguaje Natural (NLP) [4] proporciona un medio para "desbloquear" esta importante fuente de datos, convirtiendo texto no estructurado en estructurado, en datos procesables para su uso en aplicaciones de soporte a las decisiones clínicas, control de calidad y monitoreo de la salud pública. Hay que tener en cuenta que los conceptos [5] pueden ser modificados por la negación (ej, "sin temblores significativos"), por la información temporal (ej, "admitido previamente por neumonía"), por el historial familiar (ej, "antecedentes familiares de enfermedades del corazón"), o modificadores que indican que el evento en realidad no ocurrió (ej, "expuesto a tuberculosis").

Es por eso que en este trabajo utilizamos técnicas que tomen en cuenta el contexto de las narrativas más que palabras específicas, es decir métodos de aprendizaje automático como Redes Neuronales Artificiales (ANN) y Máquinas de Soporte Vectorial (SVM). El objetivo del estudio es la comparación de metodologías de clasificación para el análisis multiclase de evoluciones médicas. Para esto aplicamos las metodologías a evoluciones en las cuales identificamos si corresponden a Diabetes Tipo 1, Diabetes Tipo 2 o no se relacionan con esta afección.

Esto podría servir para levantar alarmas en el caso de que el problema asociado a la evolución no sea el correspondiente o cuando las características demarcaran una posible ocurrencia de un problema que no esté registrado en la HCE del paciente.

## **2 Antecedentes**

No sólo es importante la detección temprana de la enfermedad sino también identificar correctamente el tipo que corresponde, ya que los tratamientos pueden diferir según la sintomatología presentada.

La Diabetes Mellitus [6] es un desorden metabólico crónico caracterizado por niveles persistentemente elevados de glucosa en la sangre, como consecuencia de una alteración en la secreción y/o acción de la insulina. La Diabetes Mellitus tipo 1 (DBT1) se caracteriza por la destrucción de las células beta pancreáticas, que se traduce en un déficit absoluto de insulina y dependencia vital a la insulina exógena. La Diabetes Mellitus tipo 2 (DBT2) se caracteriza por resistencia insulínica, que habitualmente se acompaña de un déficit relativo de insulina.

Algunos trabajos analizan la detección o predicción de la enfermedad a través de variables estructuradas. Por ejemplo en el trabajo de Hera et al. [7] alimentan un modelo estadístico de regresión logística con variables como el nivel de glucemia para la detección de diabetes desconocida en pacientes coronarios. Khan et al. [8]

entrenan una red neuronal para el análisis de enfermedades renales a partir de la ocurrencia de síntomas (temperatura, sensación de náusea, dolores) para facilitar el diagnóstico del médico. En el estudio de Vassy et al. [9] utilizan variantes genómicas para la identificación de Diabetes Tipo 2 y en el de Yu et al. [10] aplican datos estructurados de contexto del paciente (como edad, peso, historial clínico) para detectar diabetes y pre-diabetes con la aplicación de técnicas de máquinas de soporte vectorial.

Sin embargo, en nuestro trabajo buscamos la detección de diabetes a partir de la información de contexto presente en la narrativa médica de las HCE. La narrativa médica siempre ha sido una parte vital de la medicina [11] y se puede entender como el puente entre la evidencia de estudios a gran escala y el arte de la medicina de la aplicación de este conocimiento al caso particular de cada paciente.

Existen varios trabajos donde se utilizan técnicas de procesamiento del lenguaje natural para la extracción de información en texto libre. Por ejemplo, en el de Chau, Xu y Chen [12] donde desarrollan un extractor basado en redes neuronales para diferenciar entidades en reportes policiales.

Pero si nos enfocamos en los trabajos con narrativas médicas, vemos por ejemplo que Breydo, Chu y Turchin [13] desarrollan un algoritmo para la detección de medicamentos inactivos para removerlos de la lista de las HCE. En Savoya et al. de Clínica Mayo [14] desarrollan un sistema llamado cTakes para Análisis de Texto y extracción del conocimiento, a partir de anotaciones semánticas de diversas estructuras de las oraciones que integran la narrativa.

Deng y Denecke [15] plantean la dificultad de los médicos de encontrar rápidamente la información del paciente siendo que la misma se encuentra distribuida en grandes cantidades de narrativa que se va incrementando a lo largo de su evolución. Para eso presentan una metodología de sumarización de la HCE a través de extracción de información, sentimental analysis, nubes de palabras y tecnologías aplicadas en creación de resúmenes.

Pero en el presente estudio no nos enfocamos en la extracción sino en la clasificación de documentos. En particular, clasificaciones de evoluciones médicas según si se relacionan con DBT Tipo 1, DBT Tipo 2 o representan casos de otras afecciones.

En el trabajo de Rodríguez, Calot y Merlino [16] el objetivo es la clasificación de texto médico en español. Los campos a clasificar son los de prescripciones médicas, que resultan campos de una oración o segmento de ésta. Para la representación de los tokens usan unigramas. Para resolver la clasificación comparan las metodologías de SVM y Naïve Bayes Multinomial.

En el trabajo de Wright et al. [17] categorizan notas de texto libre médico con SVM. Los conjuntos de entrenamiento y testeo fueron armados a partir de evoluciones médicas relacionadas o no con diabetes. La diferencia con nuestro estudio es que la clasificación que presentan es biclase, ya que entrenan la SVM para detectar Diabetes (cualquiera sea el tipo) o no.

### 3 Metodología

Nuestra metodología consistió en la comparación del entrenamiento de una Red Neuronal y una SVM multiclase para la detección y diferenciación de DBT Tipo 1, DBT Tipo 2 y Otros. Para esto formamos un conjunto de entrenamiento y uno de testeo a partir de evoluciones médicas.

En las HCE hay campos estructurados como por ejemplo los valores de laboratorio de estudios de pacientes. Pero se valora la posibilidad de que los médicos realicen anotaciones en texto libre. Uno de esos campos es de evoluciones.

Cada evolución es un texto complejo formado generalmente por varias oraciones. En cada una de ellas los médicos utilizan diferentes formas de expresión, nombran medicamentos, comportamientos de los pacientes, antecedentes propios o de familiares, frecuentemente con abreviaturas, comentarios y palabras escritas de diferentes maneras. Dependiendo de la patología y sintomatología, a veces los médicos vuelcan en las evoluciones información que ya está en otros campos. Pero muchas veces, expresan información que es valioso recuperar.

Es así que buscamos 340 ejemplos de evoluciones al azar donde el problema asociado fuera Diabetes Tipo 1 (130 casos), Diabetes Tipo 2 (130 casos) y Otros (es decir, casos que no fueran de diabetes) (80 casos).

En los casos de Diabetes, encontramos evoluciones referidas a controles de la enfermedad, cambios de medicación, complicaciones, sintomatología asociada, entre otros.

**Fig. 1.** Ejemplo de una evolución de DBT Tipo 1

Paciente de 9 años con diagnóstico de DBT desde los 2 años  
Regular control de la DBT  
Cursa 4to grado con t de aprendizaje.  
la madre refiere distención  
Va a PSG 2/sem  
Recibió ritalina a los 5 años por 2 meses Disficultades en la motricidad fina.  
  
Ef: normal  
Plan: VC informe PSG  
Sesiones TO

**Fig. 2.** Ejemplo de una evolución de DBT Tipo 2

Paciente de 63 años dbt tipo 2 de reciente dg Glucemia del 7-7-10 164 y del 20-7-10 TTOG Gluc basal 174 120'260 mg  
Hb A1c 7.5%  
Derivado por su m'édica de cabecera para adecuar plan alimentario  
Hipotiroidismo en tto con levotiroxina  
A física no realiza  
Habitualmente saltea almuerzos y concentra la mayor parte de la ingesta por sí abundante consumo de frutas  
Peso 102kg bmi 33  
Imp dg obesidad G1  
Plan alimentario hipoc con selección de H de c+ A física 150' semanas  
Continúa control con su médica de cabecera

Encontramos que los médicos no siempre asocian a la evolución el problema exacto, es decir, que vimos ejemplos donde el texto se refería a DBT T1 y el problema asociado era DBT T2 y viceversa. Esto incrementa la importancia de este trabajo, ya que una incorrecta clasificación al momento del registro del paciente puede interferir en el tratamiento adecuado.

Los casos que no son de Diabetes refieren a evoluciones que pueden deberse a fracturas, análisis oftalmológicos, controles oncológicos, etc.

**Fig. 3.** Ejemplo de una evolución de Otros

Paciente de 87 años con diagnóstico de cáncer de pulmón estadio III. En el día de ayer consulto por guardia por agudización de disnea se medicó con nebulizaciones con broncodilatadores y mejoró sintomatología. Se planificó nueva fibrobroncoscopia para el día de mañana para evaluación de stent y eventual aspirado de secreciones. Hoy concurre para realizar quimioterapia pero se pospone la misma hasta haber realizado el procedimiento previamente descrito.

Con esta selección armamos 2 corpórea con 250 evoluciones para entrenamiento (100 T1, 100 T2 y 50 otros) y 90 evoluciones para testeo (30 T1, 30 T2 y 30 otros).

### 3.1 Preprocesamiento

La representación que elegimos para modelar los datos fue la de Bag-of-Word (BOW bolsa de palabras). Para facilitar esta aproximación realizamos algunos preprocesamientos. Pasamos los textos a minúsculas (para evitar repeticiones de palabras), sacamos acentos (de la misma manera que en el caso de minúsculas, palabras acentuadas, sin acentos o mal acentuadas representan el mismo término).

También quitamos blancos extras, caracteres especiales como signos de puntuación, porcentajes (utilizados por ejemplo en la descripción de estudios) y números. Creamos una lista de "stopwords" a la cual exceptuamos conectores como "no" e "y". Esto es porque consideramos importante la información de negaciones y conjunciones, pero sacando las stopwords reducimos el tamaño de la BOW.

Otra técnica de reducción que probamos fue el "truncado" de palabras (stem en inglés), es decir dejar sólo la raíz de las mismas. Pero encontramos que esto afectaba los resultados disminuyendo la exactitud. Interpretamos que las palabras truncadas perdían información importante al entrenamiento.

### 3.2 Entrenamiento

Para la creación y entrenamiento de la Red Neuronal utilizamos el programa R [18] con la librería nnet. Para la realización de una salida multiclase codificamos los tipos de afección de la siguiente manera:

**Tabla 1.** Codificación multiclase para la red

Clase	Salida 1	Salida 2
DBT Tipo 1	1	0
DBT Tipo 2	0	1
Otros	0	0

Para reducir el tamaño de la BOW o sea, la cantidad de entradas a la red hicimos algunas pruebas preprocesando los datos con un Análisis de Componentes Principales (PCA). Esta es una transformación lineal que reduce la dimensionalidad de los datos. Al reducir el espacio bajamos notablemente el tiempo de entrenamiento, permitiendo hacer pruebas rápidas con más cantidad de neuronas en la capa oculta.

Dado que las redes neuronales dependen del tiempo de entrenamiento, comparamos los resultados con otra metodología de clasificación como lo son las Máquinas de Soporte Vectorial. En las SVM el entrenamiento es mucho menos costoso en tiempo y baja la probabilidad de sobre entrenamiento. Para la creación y entrenamiento de la SVM usamos el programa R con la librería e1071.

## 4 Resultados

Para medir los resultados del aprendizaje medimos *accuracy*, *sensitivity* y *specificity* [19] estas últimas por cada clase, ya sea Clase 0 para otros, Clase 1 para DBT Tipo 1 y Clase 2 para DBT Tipo 2.

### 4.1 Redes Neuronales

Uno de los parámetros variables fue la cantidad de neuronas de la capa intermedia de la red. A medida que probamos con corpórea más grandes, el tamaño de la BOW resultante se incrementaba por la cantidad de palabras diferentes que se pueden ver en las frases a analizar (como decíamos antes, las evoluciones están formadas muchas veces por varias oraciones o conceptos continuos). Es así que el agregado de una o varias neuronas ocultas hace que la cantidad de las conexiones de la red crezca considerablemente como así también el tiempo de entrenamiento.

Probamos entonces con diferentes cantidades de neuronas ocultas, primero comenzamos con 2 neuronas que representaron 7700 pesos sinápticos, lográndose una *accuracy* del 75,5%

**Tabla 2.** Resultados ANN con 2 n-ocultas.

Clase	Sensitivity	Specificity
General	<b>Accuracy</b>	<b>0.755</b>
Clase 0	0.7333	0.9666
Clase 1	0.7666	0.8666
Clase 2	0.7666	0.8000

Luego tratamos de mejorar el aprendizaje e incrementamos la cantidad de neuronas a 4. Esto duplicó la cantidad de pesos a un orden de 15400 y el tiempo de entrenamiento fue de unos 15 minutos aproximadamente.

**Tabla 3.** Resultados ANN con 4 n-ocultas.

Clase	Sensitivity	Specificity
General	<b>Accuracy</b>	<b>0.8555</b>
Clase 0	0.8333	0.9833
Clase 1	0.8666	0.8666
Clase 2	0.8666	0.9333

Se obtuvo una mejora considerable con una accuracy del 85,5%.

Cuando subimos la cantidad de neuronas a 5, se produjo una mejora del aprendizaje del conjunto de entrenamiento cuya accuracy llegó al 100% pero disminuyó la del testeo al 76,6%. Esto resulta en una situación de sobre entrenamiento que puede ocurrir en el caso de las redes neuronales. Con 6 neuronas intermedias, obtuvimos una accuracy del 78,8% con un tiempo de 35 minutos de aprendizaje y unos 23000 pesos sinápticos.

#### 4.2 Redes Neuronales con PCA

Para las pruebas con PCA, comenzamos con 2 neuronas intermedias. Sin PCA la cantidad de conexiones de la red era de 7700, con PCA estuvo en el orden significativamente menor de 400. Con 6 neuronas la cantidad era de 23000, con PCA registramos 1200. Entrenamos ejemplos de hasta 50 neuronas ocultas en cuestión de unos pocos minutos (con 50 neuronas, la cantidad de pesos fue del orden de 10000).

Sin embargo, la accuracy no fue tan alta como la obtenida sin PCA. Con PCA y 2 neuronas, vimos un 53,33% (contra 75,5% sin PCA), y con 4 neuronas ocultas un 63,3% (contra 85,5% obtenido sin PCA).

Valores altos de accuracy fueron encontrados recién a las 25 y 30 neuronas intermedias (cantidades más grandes de neuronas bajaban nuevamente el porcentaje).

**Tabla 4.** Resultados ANN con PCA y 25 n-ocultas.

Clase	Sensitivity	Specificity
General	<b>Accuracy</b>	<b>0.7000</b>
Clase 0	0.6666	0.8500
Clase 1	0.7666	0.8333
Clase 2	0.6666	0.8666

**Tabla 5.** Resultados ANN con PCA y 30 n-ocultas.

Clase	Sensitivity	Specificity
General	<b>Accuracy</b>	<b>0.7777</b>
Clase 0	0.7333	0.9000
Clase 1	0.8000	0.8500
Clase 2	0.8000	0.9166

Tampoco se obtuvo una mejora en la specificity (cantidad de falsos negativos) en comparación a los resultados obtenidos sin PCA.

### 4.3 Máquinas de Soporte Vectorial

Para el caso de SVM, luego de varias pruebas elegimos como una mejor respuesta un kernel sigmoide. No encontramos diferencias significativas cambiando el valor del costo, por lo que no describimos variaciones de ese parámetro en este trabajo.

Realizamos variaciones con respecto al parámetro gamma del kernel. Al ser una SVM multiclase pudimos comparar también los valores obtenidos con las técnicas de predicción. Esto es, por alternativas, donde la SVM elige la clase ganadora clasificando de a dos clases por vez. O por probabilidades donde la clase ganadora es la de mayor probabilidad obtenida [20].

**Tabla 6.** Algunos resultados de accuracy SVM con varios gamma.

Gamma	Alternativas	Probabilidades
0.1	0.6889	0.5222
0.05	0.6222	0.6222
0.01	0.5333	0.6889
0.001	0.7333	0.7333

Como en el caso de las redes neuronales, probamos preprocesar los datos con la técnica de PCA. Encontramos que los valores de accuracy se incrementaban. Los tiempos de entrenamiento no variaron significativamente porque ya eran cortos, pero podría reducir tiempos a futuro con corpórea más grandes. Los gamma que resultaron con mejores valores para SVM con PCA fueron los de 0.015 y 0.016 con probabilidades.

**Tabla 7.** Resultados SVM PCA prob gamma 0.015.

Clase	Sensitivity	Specificity
General	<b>Accuracy</b>	<b>0.8333</b>
Clase 0	0.9000	0.9500
Clase 1	0.7667	0.8833
Clase 2	0.8333	0.9167

**Tabla 8.** Resultados SVM PCA prob gamma 0.016.

Clase	Sensitivity	Specificity
General	<b>Accuracy</b>	<b>0.8444</b>
Clase 0	0.9000	0.9500
Clase 1	0.7667	0.9000
Clase 2	0.8667	0.9167

Si bien con la red neuronal sin PCA se logró una accuracy de 85,5%, la specificity de Clase 1 es levemente mejor (0.9 contra 0.86) con lo cual se podría seleccionar en el caso de buscar menos falsos positivos de DBT Tipo 1.



## 5 Conclusiones y trabajos futuros

En este trabajo presentamos una comparación de metodologías de clasificación para evoluciones médicas multiclase. En particular, clasificamos evoluciones relacionadas con Diabetes Tipo 1, Tipo 2 u otro tipo de afección. En el caso de DBT, no sólo es importante la detección temprana de la enfermedad sino también identificar correctamente el tipo que corresponde, ya que los tratamientos pueden diferir según la sintomatología presentada.

Por lo que observamos, obtuvimos una mejor accuracy con la red neuronal sin PCA, pero en poca diferencia con SVM con PCA. Si tomamos en cuenta grandes volúmenes de datos, las SVM tienen un tiempo de entrenamiento más corto (las redes pueden llegar a utilizar horas de aprendizaje con relativamente conjuntos no tan grandes). Y además, vimos que también podrían influir casos en los que se prioriza detectar menor cantidad de falsos positivos (specificity mayor).

Las SVM [21] implementan un principio de minimización del riesgo estructural en vez del riesgo empírico. Eso hace que minimicen un límite superior a la generalización de error en lugar de minimizar el error de entrenamiento, evitándose el sobre entrenamiento. Y el tiempo de aprendizaje se reduce también por el hecho de que las ANN pueden tener varios mínimos locales y las SVM sólo un mínimo global a alcanzar. Como conclusión podemos decir que las SVM con un preprocesamiento de PCA sería un clasificador útil en el caso de evoluciones médicas, que representan texto libre con anotaciones complejas en longitud, términos y formato.

En este trabajo utilizamos córpora no tan grandes a los fines de la comparación (en particular por los tiempos de la red, y la complejidad y longitud de las evoluciones). En trabajos futuros la idea es coleccionar un conjunto grande de datos para la aplicación y ajuste del método seleccionado de modo que la clasificación obtenida por esta metodología sirviera como entrada a los servicios de terminología para levantar eventuales alarmas sobre el problema asociado a cada evolución. Otros aspectos a desarrollar podrían ser el uso de otras representaciones de datos como como n-gramas y extender este estudio para otros tipos de problemas.

## Referencias

1. González Bernaldo de Quirós F., Luna D., Baum A., Plazzotta F., Otero C., Benítez S.: Incorporación de tecnologías de la información y de las comunicaciones en el Hospital Italiano de Buenos Aires. Com. Económica para América Latina y el Caribe CEPAL (2012)
2. Gambarte M.L., Lopez Osornio A., Martinez M., Reynoso G., Luna D., González Bernaldo de Quirós F.: A Practical Approach to Advanced Terminology Services in Health Information Systems. Studies in Health Technology and Informatics Volume 129: MEDINFO 2007 Pages 621 – 625 (2007)
3. Jain N.L., Friedman C.: Identification of findings suspicious for breast cancer based on natural language processing of mammogram reports. Proceedings of the AMIA Annual Fall Symposium. 1997:829-833 (1997)
4. Doan S., Conway M., Phuong T.M., Ohno-Machado L.: Natural Language Processing in Biomedicine: A Unified System Architecture Overview. Clinical Bioinformatics Methods in Molecular Biology Volume 1168, 2014, pp 275-294 (2014)

5. Friedman C., Shagina L., Lussier Y., Hripcsak G.: Automated Encoding of Clinical Documents Based on Natural Language Processing. *Journal of the American Medical Informatics Association : JAMIA*. 2004;11(5):392-402. doi:10.1197/jamia.M1552 (2004)
6. <http://www.supersalud.gob.cl/difusion/572/w3-propertyvalue-3130.html> Superintendencia de Salud Chile (consultado en julio 2015)
7. Hera J.M., Vegas J.M., Hernández E., Lozano I., García-Ruiz J.M., Fernández-Cinadevilla O.C., Carro A., Avanzas P., Torres F., Bayón J., Menéndez T., Jiménez-Navarro M., Delgado E.: Rendimiento de la glucohemoglobina y un modelo de riesgo para la detección de diabetes desconocida en pacientes coronarios. *Rev Esp Cardiol*. 2011;64:759-65. - Vol. 64 Núm.09 DOI: 10.1016 (2011)
8. Khan I.Y., Zope P.H., Suralkar S.R.: Importance of Artificial Neural Network in Medical Diagnosis disease like acute nephritis disease and heart disease. *International Journal of Engineering Science and Innovative Technology (IJESIT)* Volume 2, Issue 2, March 2013 ISSN: 2319-5967 (2013)
9. Vassy J.L., Hivert M.F., Porneala B., Dauriz M., Florez J.C., Dupuis J., Siscovick D.S., Fornage M., Rasmussen-Torvik L.J., Bouchard C., Meigs J.B.: Polygenic type 2 diabetes prediction at the limit of common variant detection. *American Diabetes Association* 2014 Jun;63(6):2172-82. doi: 10.2337/db13-1663 (2014)
10. Yu W., Liu T., Valdez R., Gwinn M., Khoury M.J.: Application of support vector machine modeling for prediction of common diseases: the case of diabetes and pre-diabetes. *BMC Medical Informatics and Decision Mak* 2010 doi: 10.1186/1472-6947-10-16 (2010)
11. Kalitzkus V., Matthiessen P.F.: Narrative-Based Medicine: Potential, Pitfalls, and Practice. *The Permanente Journal*. 2009;13(1):80-86 (2009)
12. Chau M., Xu J.J., Chen H.: Extracting Meaningful Entities from Police Narrative Reports. *Proceeding Proceedings of the 2002 annual national conference on Digital government research* Pages 1-5 Digital Government Society of North America (2002)
13. Breydo E.M., Chu J.T., Turchin A.: Identification of Inactive Medications in Narrative Medical Text. *AMIA Annual Symposium Proceedings*. 2008:66-70 (2008)
14. Savova G.K., Masanz J.J., Ogren P.V., et al. Mayo clinical Text Analysis and Knowledge Extraction System (cTAKES): architecture, component evaluation and applications. *Journal of the American Medical Informatics Association : JAMIA*. 2010;17(5):507-513. doi:10.1136/jamia.2009.001560 (2010)
15. Deng Y., Denecke K.: Summarization of EHR using information extraction, sentiment analysis and word clouds. *Deutsche Gesellschaft für Medizinische Informatik, Biometrie und Epidemiologie* doi: 10.3205/14gmds067, urn:nbn:de:0183-14gmds0672 (2014)
16. Rodríguez J.M., Calot E., Merlino H.D.: Clasificación de Prescripciones Médicas en Español. *XV Workshop de Agentes y Sistemas Inteligentes Proceedings XX Congreso Argentino de Ciencias de la Computación* ISBN 978-987-3806-05-6 (2014)
17. Wright A., McCoy A.B., Henkin S., Kale A., Sittig D.F.: Use of a support vector machine for categorizing free-text notes: assessment of accuracy across two institutions. *Journal of the American Medical Informatics Association : JAMIA*. 2013;20(5):887-890. doi:10.1136/amiajnl-2012-001576 (2013)
18. The R Project Statistical Computing <http://www.r-project.org/> (consultado en julio 2015)
19. Parikh R., Mathai A., Parikh S., Chandra Sekhar G., Thomas R. Understanding and using sensitivity, specificity and predictive values. *Indian Journal of Ophthalmology*. 2008;56(1):45-50 (2008)
20. LIBSVM: A Library for Support Vector Machines <http://www.csie.ntu.edu.tw/~cjlin/papers/libsvm.pdf> (consultado en julio 2015)
21. Lin J.Y., Cheng C.T., Chau K.W.: Using support vector machines for long-term discharge predict. *Hydrological Sciences-Journal-des Sciences Hydrologiques* Vol 51, Issue 4 (2006)